



THE UNIVERSITY *of* EDINBURGH

Edinburgh Research Explorer

Comparative transcriptomics of primary cells in vertebrates

Citation for published version:

Alam, T, Agrawal, S, Severin, J, Young, R, Andersson, R, Amer, E, Hasegawa, A, Lizio, M, Ramilowski, J, Abugessaia, I, Ishizu, Y, Noma, S, Tarui, H, Taylor, MS, Lassmann, T, Itoh, M, Kasukawa, T, Kawaji, H, Marchionni, L, Sheng, G, Forrest, ARR, Khachigian, LM, Hayashizaki, Y, Carninci, P & De Hoon, M 2020, 'Comparative transcriptomics of primary cells in vertebrates', *Genome Research*.
<https://doi.org/10.1101/gr.255679.119>

Digital Object Identifier (DOI):

[10.1101/gr.255679.119](https://doi.org/10.1101/gr.255679.119)

Link:

[Link to publication record in Edinburgh Research Explorer](#)

Document Version:

Peer reviewed version

Published In:

Genome Research

Publisher Rights Statement:

is article, published in Genome Research, is avail-able under a Creative Commons License (Attribution 4.0 Internatio

General rights

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

Take down policy

The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact openaccess@ed.ac.uk providing details, and we will remove access to the work immediately and investigate your claim.



1 Comparative transcriptomics of primary cells in vertebrates

2 Tanvir Alam¹, Saumya Agrawal², Jessica Severin², Robert S. Young^{3,4}, Robin Andersson⁵, Erik
 3 Arner², Akira Hasegawa², Marina Lizio², Jordan Ramilowski², Imad Abugessaisa², Yuri Ishizu⁶,
 4 Shohei Noma², Hiroshi Tarui⁷, Martin S. Taylor⁴, Timo Lassmann^{2,8}, Masayoshi Itoh⁹, Takeya
 5 Kasukawa², Hideya Kawaji^{2,9}, Luigi Marchionni¹⁰, Guojun Sheng¹¹, Alistair Forrest^{2,12}, Levon M.
 6 Khachigian¹³, Yoshihide Hayashizaki⁹, Piero Carninci², Michiel de Hoon²

7 ¹College of Science and Engineering, Hamad Bin Khalifa University, Doha, Qatar;

8 ²RIKEN Center for Integrative Medical Sciences, Yokohama, Japan;

9 ³Centre for Global Health Research, Usher Institute, University of Edinburgh, United Kingdom;

10 ⁴MRC Human Genetics Unit, MRC Institute of Genetics and Molecular Medicine, University of
 11 Edinburgh, United Kingdom;

12 ⁵The Bioinformatics Centre, Department of Biology & Biotech Research and Innovation Centre,
 13 University of Copenhagen, Copenhagen, Denmark;

14 ⁶RIKEN Center for Brain Science, Wako, Japan;

15 ⁷RIKEN Center for Life Science Technologies, Yokohama, Japan;

16 ⁸Telethon Kids Institute, University of Western Australia, Subiaco, Western Australia, Australia;

17 ⁹RIKEN Preventive Medicine and Diagnosis Innovation Program, Wako, Japan;

18 ¹⁰Department of Oncology, Johns Hopkins University School of Medicine, Baltimore, MD, USA;

19 ¹¹International Research Center for Medical Sciences (IRCMS), Kumamoto University,
 20 Kumamoto, Japan;

21 ¹²Harry Perkins Institute of Medical Research, and the Centre for Medical Research, University
 22 of Western Australia, QEII Medical Centre, Perth, Western Australia, Australia;

23 ¹³Vascular Biology and Translational Research, School of Medical Sciences, Faculty of
 24 Medicine, University of New South Wales, Sydney, Australia.

1 Abstract

2 Gene expression profiles in homologous tissues have been observed to be different between
3 species, which may be due to differences between species in the gene expression program in
4 each cell type, but may also reflect differences between species in cell type composition of each
5 tissue. Here, we compare expression profiles in matching primary cells in human, mouse, rat,
6 dog, and chicken using Cap Analysis Gene Expression (CAGE) and short RNA (sRNA)
7 sequencing data from FANTOM5. While we find that expression profiles of orthologous genes in
8 different species are highly correlated across cell types, many genes were differentially
9 expressed between species within homologous cell types. Expression of genes with products
10 involved in transcription, RNA processing, and transcriptional regulation was more likely to be
11 conserved, while that of genes encoding proteins involved in intercellular communication were
12 more likely to have diverged during evolution. Conservation of expression correlated positively
13 with the evolutionary age of genes, suggesting that divergence in expression levels of genes
14 critical for cell function was restricted during evolution. Motif activity analysis showed that both
15 promoters and enhancers are activated by the same transcription factors in different species. An
16 analysis of expression levels of mature miRNAs and of primary miRNAs identified by CAGE
17 revealed that evolutionary old miRNAs are more likely to have conserved expression patterns
18 than young miRNAs. We conclude that key aspects of the regulatory network are conserved,
19 while differential expression of genes involved in cell-to-cell communication may contribute
20 greatly to phenotypic differences between species.

21 Introduction

22 Vertebrate organisms consist of 100's of cell types, with more than 400 cell types defined in
23 human (Vickaryous and Hall 2006). Traditionally, cell types have been defined by their tissue of

1 origin as well as by their cellular phenotypes including morphology, staining properties, enzyme
2 histochemistry, and cell surface marker recognition by antibodies (Vickaryous and Hall 2006).
3 Cell type characterization has been supplemented by molecular approaches such as molecular
4 fingerprinting (Arendt 2008) as well as genome-wide profiling of the transcriptome of primary
5 cells (The FANTOM Consortium and the RIKEN PMI and CLST (DGT) 2014). To this end, the
6 Human Cell Atlas initiative aims to comprehensively define human cell types by performing
7 transcriptome analysis in single cells on a massive scale (Regev et al. 2017).

8

9 Evolution of anatomy is thought to primarily depend on the evolution of gene expression
10 patterns and regulation, rather than the evolution of the encoded protein sequences (King and
11 Wilson 1975; Britten and Davidson 1971). While comparative studies have shown that gene
12 expression programs in matching tissues are largely conserved between species (Chan et al.
13 2009; Brawand et al. 2011; Merkin et al. 2012; Su et al. 2002), many genes were found to be
14 differentially expressed (Su et al. 2002, Lin et al. 2014; Yue et al. 2014). Although such
15 expression differences between human and mouse for specific genes may be due in part to
16 differences in cell type composition of the analyzed tissues (Breschi et al. 2017), little overlap
17 was found in terms of differentially expressed genes between human and mouse in dynamic
18 studies of primary cells during erythropoiesis (Pishesha et al. 2014), and of primary
19 macrophages upon stimulation by lipopolysaccharide (Schroder et al. 2012) or by glucocorticoid
20 (Jubb et al. 2016). Collectively, these findings suggest that also in matching primary cells many
21 genes are differentially expressed between species. As cells with an identical cellular phenotype
22 may display distinct and disparate molecular phenotypes, the question what key transcriptomic
23 features define a cell type is raised (Arendt et al. 2016).

24

25 The confounding effects of cell type composition in tissue-based studies can be avoided by
26 comparing the transcriptome of different species in homologous primary cells. Here, we present

1

1 a comparative analysis of genome-wide expression in vertebrate species profiled in FANTOM5
2 (The FANTOM Consortium and the RIKEN PMI and CLST (DGT) 2014; Lizio, Mukarram, et al.
3 2017; Lizio, Deviatiiarov, et al. 2017) to elucidate patterns of gene expression conservation
4 during evolution.

5 Results

6 The FANTOM5 collection contains Cap Analysis Gene Expression (CAGE) data for 3 primary
7 cell types in human, mouse, rat, dog, and chicken, and for an additional 12 cell types in human
8 and mouse only (Supplemental Table S1). We identified 15,538, 14,915, 13,759, and 8,696
9 protein-coding genes in mouse, rat, dog, and chicken, respectively, with a one-to-one
10 orthologous gene in human, and 6,561 protein-coding genes with one-to-one orthologs in all five
11 species (see Methods for details). Principal Component Analysis (PCA) of all human and mouse
12 samples revealed a liver-specific cluster, a mesenchymal cluster, and a hematopoietic cluster
13 (Fig. 1A), and similarly PCA for cell types with CAGE data available in all five species showed a
14 hepatocyte cluster and a mesenchymal cluster (Fig. 1B). Within each cluster, samples tended to
15 cluster by species (Fig. 1), consistent with the "species signal" phenomenon observed
16 previously (Musser & Wagner 2015).

1 Expression levels of pairs of orthologous genes were positively correlated across cell types, with
2 median Pearson's correlation values ranging from 0.38 to 0.72 ($p < 10^{-100}$, mouse, rat, and dog;
3 $p = 2.2 \times 10^{-42}$, chicken; Fig. 2A,B). Nevertheless, in specific cell types we found significant
4 differences in absolute expression of orthologs in different species (Fig. 2A,C). Pairwise
5 differential expression analysis between genes in human and their orthologs in mouse, rat, dog,
6 or chicken for each primary cell type in FANTOM5 revealed that on average 52% of expressed
7 genes were differentially expressed (Benjamini-Hochberg corrected $p < 0.1$) between the two
8 species (Fig. 2C and Supplemental Table S2).

9

10 In each species, we defined the dominant promoter for each gene as the most highly expressed
11 promoter associated with the gene. The genomic region of the dominant promoter of more than
12 80% of genes in mouse, rat, and dog, and 50% of genes in chicken had an orthologous region
13 in the human genome; the majority of those overlapped the corresponding human dominant
14 promoter (Fig. 3A). Genes were more likely to be differentially expressed if their dominant
15 promoter was located in a genomic region that did not have an orthologous genome sequence
16 in the human genome (Fisher combined $p < 10^{-100}$; Fig. 3B and Supplemental Fig. S1),
17 suggesting that gain or loss of promoter sequence regions during evolution contributes to the
18 emergence of gene expression differences between species.

19

20 We hypothesized that genes critical for cellular functioning would both be more conserved and
21 have more conserved expression patterns, and indeed found the expression levels of
22 evolutionarily older genes to be more conserved (Fisher combined $p < 10^{-100}$; Fig. 4A and
23 Supplemental Fig. S2). Gene Ontology analysis of differentially expressed genes showed that
24 genes with products involved in transcription, RNA processing, and transcriptional regulation,
25 were more likely to have conserved expression levels, whereas genes encoding proteins
26 localized to the plasma membrane and extracellular space as well as signaling proteins were

1

1 most likely to be differentially expressed (Fig. 4B and Supplemental Table S3). This suggests
2 that the transcriptional program in each cell tends to be conserved during evolution, while genes
3 in the periphery of the transcriptional regulatory network, especially those involved in cellular
4 communication, tend to diverge in expression.

5

6 As an independent confirmation, we applied integrative correlation analysis (Parmigiani et al.
7 2004) by first calculating the correlations across cell types between all genes for human and
8 mouse separately, and then the correlation across orthologous genes between corresponding
9 rows in these two correlation matrices. This yielded the correlation-of-correlations, or integrative
10 correlation coefficient, as a measure of the degree of expression conservation during evolution
11 for each gene. We then ranked genes based on their integrative correlation coefficient, and
12 performed gene set enrichment analysis to identify biological processes most conserved or
13 most divergent between the two species (see Methods section).

14

15 The integrative correlation coefficient values ranged between -0.52 to 0.59, and their observed
16 distribution was skewed to the right, with a median of 0.25 (Supplemental Fig. S3A and
17 Supplemental Table S4), suggesting that, overall, gene expression profiles tend to be
18 conserved between human and mouse. Similarly to our conclusions for Gene Ontology analysis
19 of differentially expressed genes, fundamental cellular processes involved in cell homeostasis
20 and maintenance tended to rank higher in integrative correlation analysis, while gene sets
21 encompassing processes associated with cell-to-cell signaling and other biological processes
22 taking place in the extra-cellular space (e.g., neuronal and synapse development) were more
23 likely to rank lower, suggesting less conserved underlying networks (Supplemental Fig. S3B and
24 Supplemental Table S4).

25

1 As a complement to the differential gene expression analysis, we calculated the correlation
2 across genes for each cell type and species. Expression levels were positively correlated within
3 each species as well as between species for related cell types (Supplemental Fig. S4, S5),
4 suggesting the relative ranking of genes by their expression level tends to be conserved. The
5 correlation value decreased exponentially as a function of phylogenetic distance between
6 species, and dropped off most rapidly for mesenchymal stem cells compared to aortic smooth
7 muscle cells and hepatocytes (Supplemental Fig. S6). Consistent with the differential gene
8 expression results, expression levels were more highly correlated for genes for which the
9 dominant promoter had an orthologous genome region in human compared to genes for which
10 the dominant promoter did not have an orthologous genome region (Fisher combined $p < 10^{-100}$;
11 Supplemental Fig. S7 and S8), as well as for evolutionarily ancient genes compared to recent
12 genes (Fisher combined $p < 10^{-100}$; Supplemental Fig. S9 and S10). A Gene Ontology analysis
13 of correlation values again showed that genes with functional roles associated with RNA biology
14 in the nucleus tended to have conserved expression levels, while genes with functions
15 associated with the plasma membrane, extracellular space, and signaling had lower correlation
16 values (Supplemental Fig. S11 and Supplemental Table S3).

17

18 To confirm these findings in an independent gene expression data set, we performed differential
19 expression analysis on previously published RNA-seq expression data for endometrial stromal
20 fibroblast primary cells in human, rat, rabbit, ferret, cow, and opossum (Kin et al. 2016). We
21 again found that evolutionarily ancient genes were more likely to have conserved expression
22 levels compared to recent genes (Fisher combined $p = 1.0 \times 10^{-11}$; Supplemental Fig. S12).
23 Results of Gene Ontology analysis of differentially expressed genes for these data were highly
24 consistent with those observed in the FANTOM5 samples (Supplemental Fig. S12), including
25 evidence of rapid evolution of signaling pathways as observed previously (Kin et al., 2016). A
26 comparative analysis of RNA-seq expression data in matching tissues in human and mouse

(The ENCODE Project Consortium 2012) also showed preferential conservation of expression levels of evolutionarily ancient genes (Supplemental Fig. S13A,B), and yielded similar patterns of Gene Ontology enrichment (Supplemental Fig. S13C).

To understand how evolution of the transcriptional regulatory network affects evolution of gene expression, we used the MotEvo sequence motif analysis software (Arnold et al. 2012) for the 190 motifs compiled in SwissRegulon (Pachkov et al. 2013) to identify potential transcription factor binding sites (TFBSs) in the human, mouse, rat, dog, and chicken genomes. We evaluated the TFBS prediction accuracy using ChIP-seq data (Supplemental Table S5) for transcription factors associated with each motif (Supplemental Fig. S14). Conservation between species of the expression patterns of orthologous genes depended on the concordance in TFBS presence in the promoter of each gene (Supplemental Fig. S15), demonstrating the contribution of cis-regulatory evolution to expression divergence between species. To analyze trans-regulatory evolution, we performed motif activity analysis (Suzuki et al. 2009), which uses linear decomposition of genome-wide gene expression patterns based on the TFBSs found in the promoter of each gene, resulting in motif activities representing the average expression level of genes with a predicted binding site for each motif. Fig. 5 shows the broadly expressed transcription factor TP53 (Fig. 5A), the hematopoietic lineage-specific RUNX transcription factors (Fig. 5B), and the motif associated with the hepatocyte-specific HNF4A transcription factor (Fig. 5C) as examples of motifs with activities highly correlated between human and mouse. In contrast, the motif associated with the testis-specific transcription factor SPZ1 did not show evidence of activation either in human or mouse, as testis was not included in our samples (Fig. 5D). In general, motif activities were highly correlated across samples between human and mouse ($p = 5.5 \times 10^{-25}$, Mann-Whitney U test), rat ($p = 3.9 \times 10^{-9}$), dog ($p = 4.5 \times 10^{-6}$), and chicken ($p = 9.2 \times 10^{-4}$), compared to randomized pairs of motifs (Fig. 5E and Supplemental Table S6).

1

1

2 We then asked if enhancers likewise were activated by the same transcription factors in
3 different species. Enhancers were previously identified in human and mouse from FANTOM5
4 CAGE data by searching for a characteristic bidirectional expression pattern (Andersson et al.
5 2014). We predicted enhancers in rat, dog, and chicken by applying the same pipeline on the
6 FANTOM5 CAGE data in these species (Supplemental Table S7), and used the CAGE
7 expression level at each enhancer as a measure of its activity (Andersson et al. 2014). For each
8 species, the motif activity calculated from gene promoter expression profiles correlated with the
9 motif activity based on enhancer expression profiles (human, $p = 1.2 \times 10^{-20}$, Mann-Whitney U
10 test), mouse ($p = 5.6 \times 10^{-22}$), rat ($p = 5.6 \times 10^{-5}$), dog ($p = 2.5 \times 10^{-5}$), and chicken ($p = 5.7 \times$
11 10^{-4}) (Supplemental Fig. S16), indicating that in each species enhancers are activated by the
12 same transcription factors as promoters. Between species, the motif activity calculated from
13 enhancer expression profiles were highly correlated between human and mouse ($p = 1.6 \times$
14 10^{-18} , Mann-Whitney U test), rat ($p = 2.6 \times 10^{-6}$), dog ($p = 0.0032$), and chicken ($p = 0.044$) (Fig.
15 5F and Supplemental Table S6). We conclude that both promoters and enhancers are activated
16 by the same transcription factors in different species.

17

18 Next, we extended our comparative analysis to the expression levels of microRNAs (miRNAs).
19 miRNAs are small non-coding RNA (typically 22 nts) that silence mRNA post-transcriptionally
20 and regulate biological processes such as cell growth and differentiation by functional effects on
21 direct targets and regulatory networks (Bracken et al. 2016). In the FANTOM5 collection, short
22 RNA (sRNA) sequencing data for matching primary cell types in different species were available
23 for aortic smooth muscle cells (Supplemental Table S1 and Supplemental Table S8). We
24 annotated known (Supplemental Table S9) and candidate novel (Supplemental Table S10)
25 miRNAs in rat, dog, and chicken in the same way as done previously (De Rie et al. 2017) for
26 human and mouse. Differential expression analysis between human and mouse, rat, dog or

1 chicken showed that about half of the orthologous miRNAs had statistically significant different
2 expression levels in the two species (Fig. 6A and Supplemental Table S11). Dividing miRNAs
3 into three categories based on their evolutionary age revealed that evolutionarily older miRNAs
4 were more likely to have conserved expression levels than younger miRNAs (Fisher combined p
5 = 1.2×10^{-4}) (Fig. 6C).

6

7 Previously, we showed that CAGE data can be used to reliably infer the promoter of the primary
8 miRNA (pri-miRNA) transcript, and that the corresponding CAGE expression levels can be used
9 as a proxy for the expression level of the mature miRNA (De Rie et al. 2017). We manually
10 curated pri-miRNA promoters previously identified computationally for mouse (De Rie et al.
11 2017), and, using the same approach, identified pri-miRNA promoters for miRNAs in rat, dog,
12 and chicken (Supplemental Table S12). In aortic smooth muscle cells, expression levels of the
13 mature miRNA measured by sRNA sequencing correlated with the CAGE expression level of
14 the pri-miRNA for mouse, rat, dog, and chicken (Supplemental Fig. S17). The curated primary
15 miRNA promoter annotations as well as expression levels of the mature and primary miRNA are
16 visualized and available for download through an interactive web interface at
17 https://fantom.gsc.riken.jp/zenbu/reports/#FANTOM_miRNA_atlas.

18

19 Using these promoters together with previously curated pri-miRNA promoters for human (De
20 Rie et al. 2017), we performed differential expression analysis of miRNAs in human compared
21 to mouse, rat, dog, and chicken. In aortic smooth muscle cells in mouse, rat, dog, and chicken,
22 log-ratios of mature miRNA expression levels, as measured by sRNA sequencing, correlated
23 well with the log-ratios for pri-miRNAs, as measured by CAGE expression data (Supplemental
24 Fig. S18), and among the miRNAs differentially expressed in both data sets, more than 80%
25 showed concordant up- or down-regulation of the mature miRNA and the pri-miRNA, suggesting

1

1 that few of the identified differentially expressed miRNAs were false positives (Supplemental
2 Fig. S18).

3

4 Differential CAGE expression analysis of pri-miRNAs revealed that the majority of expressed
5 orthologous miRNAs have different expression levels in human compared to mouse, rat, dog,
6 and chicken (Fig. 6B and Supplemental Table S13), consistent with the results obtained for
7 mature miRNAs (Fig. 6A). We found significantly fewer differentially expressed miRNAs for
8 evolutionarily old miRNAs compared to evolutionarily recent miRNAs for 12 out of 24 pairwise
9 comparisons, a further 7 showed the same pattern without reaching statistical significance, 5
10 showed an opposite pattern without reaching statistical significance, and 0 showed a statistically
11 significant opposite pattern (Fisher combined $p = 4 \times 10^{-12}$; Fig. 6D). Therefore, using CAGE as
12 a proxy for miRNA expression allowed us to demonstrate that the patterns observed for mature
13 miRNAs by sRNA sequencing for a single cell type (Fig. 6C) can be found across a wide variety
14 of cell types.

1 Discussion

2 Comparative studies have shown considerable differences in the gene expression levels in
3 matching tissues of different species (Su et al. 2002; Lin et al. 2014; Yue et al. 2014), which at
4 least in part is due to differences in tissue composition between species (Breschi et al. 2017).
5 However, our analysis reveals that this cannot be the sole explanation, as considerable
6 expression level differences are also observed between matching primary cell types, indicating
7 that the same cellular phenotype associated with traditionally defined cell types can be achieved
8 by widely different molecular networks.

9

10 Our findings suggest that expression levels of regulators tend to be conserved across species,
11 while genes peripheral in the regulatory network, especially those involved in cellular
12 communication, are more likely to have divergent expression patterns. Previously reported
13 examples include the terminal differentiation of erythroid precursors from early to late
14 erythroblasts, where the same transcriptional regulators and other proteins important for
15 erythropoiesis were induced or repressed in human and mouse, suggesting that the core
16 regulatory program of erythroid differentiation remained conserved (Pishesha et al. 2014). In
17 contrast, genes regulated during development showed a different response between human and
18 mouse (Pishesha et al. 2014), indicating that the response of genes to the regulators of
19 erythropoiesis had evolved since the evolutionary split of human and mouse. Similarly,
20 comparing lipopolysaccharide-stimulated macrophages between human and mouse showed
21 enriched differences in the transcriptome of genes encoding proteins involved in cellular
22 communication such as cell surface receptors, inflammatory cytokines, chemokines, and their
23 intracellular signaling pathways (Schroder et al. 2012). Phenotypic differences between species

1

1 at the organismal level may thus be primarily due to differences in the interaction between cells
2 (Ramilowski et al. 2015).

3

4 Orthologous transcription factors typically recognize the same DNA sequence motif in human
5 and mouse (Cheng et al. 2014), as changes in the consensus motif during evolution would
6 simultaneously affect a large number of genes and may be too disruptive. By the same
7 argument, we can expect that expression levels of transcriptional regulators to be conserved
8 between species. As a salient example of the conservation of regulatory programs, we
9 previously found that human enhancer sequences could be activated by orthologous
10 transcription factors in corresponding tissues in human and zebrafish (Andersson et al. 2014). In
11 contrast, genomic binding sites of conserved transcription factors have diverged extensively
12 between human and mouse (Odom et al. 2007), suggesting a rewiring of the peripheral
13 regulatory network during evolution.

14

15 Due to their modular nature, enhancer regulatory elements are particularly amenable to
16 rewiring, as their cell type- and state-specific usage (Andersson et al. 2014) allow changes in
17 their regulatory connections in specific conditions while avoiding pleiotropic deleterious effects
18 on the organism in general (Carroll 2008). For example, differences in the transcriptome
19 response of human and mouse primary macrophages stimulated by glucocorticoid were
20 previously found to be associated with the turnover of glucocorticoid receptor binding sites at
21 enhancers (Jubb et al. 2016). Similarly, the cell type- and state-specific usage of different
22 promoters associated with a gene (The FANTOM Consortium and the RIKEN PMI and CLST
23 (DGT) 2014) avoids the constraints placed by pleiotropy and allows gain and loss of promoters
24 to contribute significantly to the evolution of species (Young et al. 2015). In our analysis, we
25 indeed find that changes in gene expression levels are associated with the gain and loss of
26 promoter sequence regions during evolution.

1

1

2 Our analysis further shows that the conservation of regulatory programs is not limited to
3 transcriptional regulation but extends to miRNAs. Our comparative analysis of miRNA
4 expression revealed that older miRNAs are more likely to have conserved expression levels
5 than more recent miRNAs, suggesting that highly conserved miRNAs have stronger
6 evolutionary constraints on their expression levels. As an example, we found conservation
7 across human, mouse, rat, dog, and chicken of pri-miRNA expression levels in aortic smooth
8 muscle cells of miR-22, which modulates a range of target genes including *MECP2*, *HDAC4* and
9 *MECOM* and is a key regulator of smooth muscle cell phenotype switching and neointima
10 formation (Yang et al. 2018).

11

12 The Human Cell Atlas aims to create a comprehensive map of cell types in the human body by
13 profiling gene expression levels in single cells from healthy tissues (Regev et al. 2017). Our
14 comparative analysis suggests that differences in the regulatory signature (Arendt et al. 2016),
15 rather than the overall gene expression patterns, are the key requirement for distinguishing cell
16 types.

17 Methods

18 Genome assembly version

19 For consistency with previous FANTOM5 publications (The FANTOM Consortium and the
20 RIKEN PMI and CLST (DGT) 2014; Lizio, Mukarram, et al. 2017; Lizio, Deviatiiarov, et al.
21 2017), we used genome assemblies hg19 (human), mm9 (mouse), rn6 (rat), canFam3 (dog),
22 and galGal5 (chicken) for our analysis. Previously it was shown that 99.75% and 99.94% of
23 CAGE peaks in human and mouse, respectively, could be converted unambiguously to the

1 recent genome versions hg38 (human) and mm10 (mouse), with an expression correlation
 2 value larger than 0.99 both for human and mouse (Abugessaisa et al., 2017). For 231 (human)
 3 and 202 (mouse) miRNAs included in the comparative analysis shown in Supplemental Table
 4 S13, the genomic distance between each pre-miRNA and the corresponding pri-miRNA
 5 promoter (Supplemental Table S12) was identical between genome assembly versions for 216
 6 (human) and 199 (mouse) miRNAs, and differed by less than 10 base pairs for 227 (human)
 7 and 202 (mouse) miRNAs, suggesting that the genome assembly version used had minimal
 8 effects on the analysis results.

9 Identification of orthologous genes

10 For each gene in mouse, rat, dog, and chicken defined in Ensembl (Zerbino et al. 2018) release
 11 85, we retrieved the orthologous human gene, if defined, in an 'ortholog_one2one' relationship
 12 with it in the Ensembl Compara multi-species database (Vilella et al. 2009). This yielded 16,217
 13 (human-mouse), 15,486 (human-rat), 15,861 (human-dog), 11,950 (human-chicken), and
 14 10,237 (in all five species) pairs of orthologous genes, of which 15,893 (human-mouse), 15,207
 15 (human-rat), 15,482 (human-dog), 11,873 (human-chicken), and 10,208 (in all five species)
 16 were protein-coding. Using the most recent Ensembl release available for each genome
 17 assembly (release 75 for human genome assembly hg19, release 67 for mouse genome
 18 assembly mm9, release 85 for rat genome assembly rn6 and dog genome assembly canFam3,
 19 and release 92 for chicken genome assembly galGal5), we obtained the transcription start site
 20 for all transcripts associated with each gene, defined a ± 500 base pair promoter region around
 21 each transcription start site, and merged overlapping regions. Genes for which any of the
 22 associated regions had more than 10% unidentified nucleotides (N) in their genome sequence
 23 were removed from the analysis. The number of remaining orthologous protein-coding genes

1

1 was 15,538 (human-mouse), 14,915 (human-rat), 13,759 (human-dog), 8,696 (human-chicken),
2 and 6,561 (in all five species).

3 Gene expression analysis

4 Gene expression quantitation is described in detail in the Supplemental Methods. Differential
5 gene expression analysis was performed on the raw counts using DESeq2 (Love et al. 2014)
6 version 1.22.1 with a threshold of 0.1 on the Benjamini-Hochberg adjusted p -value. PCA as well
7 as all correlation calculations (except in Integrative Correlation Coefficient analysis, described
8 below) were performed on variance-stabilized gene expression data generated as follows. First,
9 we used DESeq2 (Love et al. 2014) version 1.22.1 to estimate, for each cell type in each
10 species, the asymptotic dispersion of expression counts between replicates, and then calculated
11 its average value α across cell types and species. Next, we calculated the total tag count for
12 each sample, divided these totals by their median across samples to obtain the normalization
13 factors, and divided the counts of each sample by the corresponding factor to obtain normalized
14 count data x . We then applied the variance-stabilizing transformation (Anders & Huber, 2010) to
15 the normalized count data x :

$$16 \quad x' = \frac{2 \operatorname{arcsinh}(\sqrt{\alpha x}) - \log \alpha - \log 4}{\log 2}$$

17 The variance-stabilized gene expression data x' were averaged across replicates for each cell
18 type and for each species.

19 For each pairwise comparison in Fig. 2B, we calculated Pearson's correlation across cell types
20 between each pair of orthologous genes. Next, we randomly permuted the gene pairings,
21 calculated the correlation across cell types to find the background distribution, and performed
22 the Mann-Whitney U test comparing the set of correlation values for pairs of orthologous genes
23 to the set of correlation values for randomly permuted pairs. We also calculated a background

1

1 distribution for pairs of orthologous genes after permuting the samples, as well as the
2 cumulative distribution of correlation values for an uncorrelated bivariate normal distribution.
3 For the pairwise comparisons shown in Supplemental Fig. S4–S11, we calculated Pearson's
4 correlation between the two species for each cell type across orthologous genes. For
5 Supplemental Fig. S15, we calculated Pearson's correlation between the two species for each
6 pair of orthologous genes across cell types.

7 Promoter conservation analysis

8 Orthologous genomic regions of promoters across species were identified by applying liftOver
9 (Hinrichs et al. 2006) on chain files downloaded from the University of California, Santa Cruz
10 website (<http://hgdownload.cse.ucsc.edu/downloads.html>).

11 Gene conservation analysis

12 For each gene in human, we identified the HomoloGene group of homologous genes to which it
13 belonged in release 68 of the NCBI HomoloGene database (NCBI Resource Coordinators
14 2018). If the HomoloGene group included mammals only or vertebrates only, then the gene was
15 classified as restricted to mammals or restricted to vertebrates, respectively. Alternatively, the
16 gene was classified as conserved in bilateria if the HomoloGene group included bilateria in non-
17 vertebrate lineages. To assess the statistical significance of the increase or decrease in
18 conservation of expression in the three classes, the bilaterian, vertebrate, and mammalian class
19 were represented by an equidistant indicator variable, and applied the maximum likelihood
20 method to fit a linear regression model under the Poisson distribution to the number of
21 differentially expressed genes in each class. The corresponding p -value was calculated using
22 the likelihood-ratio test. The overall p -value was calculated by combining the p -values for the
23 pairwise comparisons using Fisher's method.

1 Gene Ontology analysis

2 Gene Ontology annotations were downloaded on June 10, 2018 from the GOA database
3 (Huntley et al. 2015). Statistical significance of over- or underrepresentation of a Gene Ontology
4 term among differentially expressed genes was calculated using Fisher's exact test, where an
5 expression-matched background was created by selecting the 10 closest genes in expression in
6 human for each differentially expressed gene. The overall p -value was calculated by combining
7 the p -values for the pairwise comparisons using Fisher's method.

8 RNA-seq expression data analysis

9 Accession numbers for ENCODE (The ENCODE Project Consortium 2012) and endometrial
10 stromal fibroblast (Kin et al. 2016) RNA-seq gene expression data are provided in the
11 Supplemental Methods. Gene conservation and Gene Ontology analysis of these data sets
12 were performed as described above.

13 Integrative Correlation Coefficient analysis

14 Integrative Correlation Coefficient analysis (Parmigiani et al. 2004) ranks genes based on the
15 degree to which their expression profiles are comparable between datasets.

16 For human and mouse separately, we constructed a CAGE expression matrix (normalized to
17 t.p.m.) for the 15,538 genes in common between human and mouse, averaging biological
18 replicates by taking the median, and performed quantile normalization separately for each
19 expression matrix. Next, we calculated the correlation between each pair of genes, again for
20 human and mouse separately, across cell types to obtain one correlation matrix for human and
21 one correlation matrix for mouse. We then calculated Pearson's correlation between human and
22 mouse for corresponding rows in these two correlation matrices to obtain the correlation-of-

1 correlations, or integrative correlation coefficient, for each gene. The null distribution was
2 obtained by randomly permuting samples 10,000 times, as described previously (Parmigiani et
3 al. 2004), using MergeMaid (Cope et al. 2004) version 2.56.0. Analysis of Functional Annotation
4 (AFA) (Kortenhorst et al. 2013, Marchionni et al. 2017, Ross et al. 2011) was conducted by
5 performing a one-sided Wilcoxon rank-sum test to compare the integrative correlation coefficient
6 values of genes in each cellular component (CC) and biological process (BP) Gene Ontology
7 category (extracted using the “org.Hs.eg.db” R/Bioconductor package version 3.8.2), requiring
8 at least 10 genes, to those of remaining genes, using the Benjamini-Hochberg multiple testing
9 correction method. All analyses were performed using the R/Bioconductor “RTopper” package
10 (version 1.30.0) (Tyekucheva et al. 2011).

11 Multiple genome alignment, TFBS prediction, and motif activity 12 analysis

13 The 100-way multiple genome alignment of human genome assembly hg19 against 99
14 vertebrate species and the 30-way multiple genome alignment of the mouse genome assembly
15 mm9 against 29 vertebrate species were downloaded from the University of California, Santa
16 Cruz website (<http://hgdownload.cse.ucsc.edu/downloads.html>), the species in the 30-way
17 mouse alignment being a subset of the species in the 100-way human alignment. For the same
18 set of 30 species, we performed pairwise genome alignments of the rat, dog, and chicken
19 genome against each of the 29 remaining species for the genome assemblies listed in
20 Supplemental Table S14 (see Supplemental Methods for details). Pairwise alignments were
21 merged into a multiple genome alignment using multiz (Blanchette et al. 2004) version 11.2
22 using the phylogenetic tree of the 30 species was extracted from the 191-way phylogenetic tree
23 in 191way.nh distributed as part of the UCSC Genome Browser bioinformatics utilities (Kuhn et
24 al. 2013) release 366 (June 5, 2018). Genome-wide TFBS predictions and motif activity analysis

1

1 were performed as described previously (Arner et al. 2015), with minor modifications as
2 described in the Supplemental Methods. The multiple genome alignment files, genome-wide
3 locations and scores of predicted TFBSs, and motif activity scripts are available at
4 http://fantom.gsc.riken.jp/5/suppl/Alam_et_al_2020/; motif activity scripts are also included in the
5 Supplemental Code.

6 Enhancer identification

7 The previously calculated set of permissive enhancers (Arner et al. 2015) was used for human
8 (65,423 enhancers) and mouse (44,459 enhancers). For rat, dog, and chicken, we first created
9 a mask for all ± 500 base pair windows around the 5' end of transcripts in the NCBI Entrez Gene
10 database (Brown et al. 2015), downloaded on 13 November 2017, as well as all windows within
11 200 base pairs of exons defined in the same database. We then applied the `bidir_enhancers`
12 script (Andersson et al. 2014) to all FANTOM5 CAGE libraries in rat, dog (Lizio, Mukarram, et
13 al. 2017), and chicken (Lizio, Deviatiiarov, et al. 2017) using the calculated mask, resulting in
14 9,372 (rat), 10,649 (dog), and 44,625 (chicken) enhancers.

15 MicroRNA analysis

16 Short RNA libraries were produced, sequenced, and processed as described previously (De Rie
17 et al. 2017) using the same RNA samples as used for CAGE expression profiling (Lizio,
18 Deviatiiarov, et al. 2017; Lizio, Mukarram, et al. 2017). Short RNA libraries not described
19 previously are listed with their matching CAGE library in Supplemental Table S1. Annotation of
20 miRNAs, candidate novel miRNA prediction, and miRNA promoter identification were performed
21 as described in the Supplemental Methods. Orthologous miRNAs were identified by performing
22 global alignment of mature miRNA sequences between species, followed by manual curation.

1

1 The evolutionary age of miRNAs was established based on the set of species in which miRNAs
2 of each family were annotated in miRBase release 21 (Kozomara and Griffiths-Jones 2014).

3 Data access

4 All raw and processed sequencing data generated in this study have been submitted to the DNA
5 Data Bank of Japan (DDBJ; <https://www.ddbj.nig.ac.jp/>) under accession number DRA008211.
6 All custom scripts generated in this study are available as Supplemental Code.

7 Competing interest statement

8 The authors declare no competing interests.

9 Acknowledgements

10 This work was supported by the following grants: Research Grant for RIKEN Omics Science
11 Center from MEXT to Y.H.; Grant of the Innovative Cell Biology by Innovative Technology (Cell
12 Innovation Program) from the MEXT to Y.H.; Research Grant from MEXT to the RIKEN Center
13 for Life Science Technologies; Research Grant from MEXT to the RIKEN Center for Integrative
14 Medical Sciences; Research Grant to RIKEN Preventive Medicine and Diagnosis Innovation
15 Program from MEXT to Y.H.; NIH-NCI award R01CA200859 and the JSPS Fellowship S19058
16 to L.M. We gratefully acknowledge the computational resources of the HOKUSAI
17 supercomputer system provided by RIKEN under project number Q18305/Q19305, which
18 enabled us to perform the pairwise genome alignments as well as the genome-wide TFBS
19 predictions by MotEvo.

20 *Author contributions:*

1 T.A., L.M., and M.d.H. analyzed the data with the help of J.R., R.A., E.A., R.S.Y., M.S.T., T.L.,
2 A.H., H.K.; G.S. provided RNA samples; Y.I., S.N., H.T., and M.I. produced the sequencing
3 libraries; I.A., M.L., H.K., and T.K. managed the data; S.A. curated the primary miRNA
4 annotations; J.S. created the interactive web interface for miRNA expression visualization;
5 A.R.R.F. and M.d.H. designed the study; T.A., L.M., and M.d.H. wrote the manuscript with the
6 help of L.M.K. and P.C.; P.C. and Y.H. supervised the FANTOM5 project.

1 References

- 2 Abugessaisa I, Noguchi S, Hasegawa A, Harshbarger J, Kondo A, Lizio M, Severin J, Carninci
3 P, Kawaji H, Kasukawa T. 2017. FANTOM5 CAGE Profiles of Human and Mouse Reprocessed
4 for GRCh38 and GRCm38 Genome Assemblies. *Scientific Data* **4**: 170107.
- 5 Andersson R, Gebhard C, Miguel-Escalada I, Hoof I, Bornholdt J, Boyd M, Chen Y, Zhao X,
6 Schmidl C, Suzuki T, et al. 2014. An atlas of active enhancers across human cell types and
7 tissues. *Nature* **507**: 455–461.
- 8 Arendt D. 2008. The evolution of cell types in animals: emerging principles from molecular
9 studies. *Nat Rev Genet* **9**: 868–882.
- 10 Arendt D, Musser JM, Baker CVH, Bergman A, Cepko C, Erwin DH, Pavlicev M, Schlosser G,
11 Widder S, Laubichler MD, Wagner GP. 2016. The origin and evolution of cell types. *Nat Rev*
12 *Genet* **17**: 744–757.
- 13 Arner E, Daub CO, Vitting-Seerup K, Andersson R, Lilje B, Drabløs F, Lennartsson A,
14 Rönnerblad M, Hrydziuszko O, Vitezic M, et al. 2015. Transcribed enhancers lead waves of
15 coordinated transcription in transitioning mammalian cells. *Science* **347**: 1010–1014.
- 16 Arnold P, Erb I, Pachkov M, Molina N, Van Nimwegen, E. 2012. MotEvo: integrated Bayesian
17 probabilistic methods for inferring regulatory sites and motifs on multiple alignments of DNA
18 sequences. *Bioinformatics* **28**: 487–494.
- 19 Blanchette M, Kent WJ, Riemer C, Elnitski L, Smit AFA, Roskin KM, Baertsch R, Rosenbloom
20 K, Clawson H, Green ED, et al. 2004. Aligning multiple genomic sequences with the threaded
21 blockset aligner. *Genome Res* **14**: 708–715.
- 22 Bracken CP, Scott HS, Goodall GJ. 2016. A network-biology perspective of microRNA function

1 and dysfunction in cancer. *Nat Rev Genet* **17**: 719–732.

2 Brawand D, Soumillon M, Necsulea A, Julien P, Csárdi G, Harrigan P, Weier M, Liechti A,
3 Aximu-Petri A, Kircher M, et al. 2011. The evolution of gene expression levels in mammalian
4 organs. *Nature* **478**: 343–348.

5 Breschi A, Gingeras TR, Guigó, R. 2017. Comparative transcriptomics in human and mouse.
6 *Nat Rev Genet* **18**: 425–440.

7 Britten RJ, Davidson EH. 1971. Repetitive and non-repetitive DNA sequences and a speculation
8 on the origins of evolutionary novelty. *Q Rev Biol* **46**: 111–138.

9 Brown GR, Hem V, Katz KS, et al. 2015. Gene: a gene-centered information resource at NCBI.
10 *Nucleic Acids Res* **43**(Database issue): D36–D42.

11 Carroll SB. 2008. Evo-devo and an expanding evolutionary synthesis: a genetic theory of
12 morphological evolution. *Cell* **134**: 25–36.

13 Chan ET, Quon GT, Chua G, Babak T, Trochesset M, Zirngibl RA, Aubin J, Ratcliffe MJ, Wilde
14 A, Brudno M, et al. 2009. Conservation of core gene expression in vertebrate tissues. *J Biol* **8**:
15 33.

16 Cheng Y, Ma Z, Kim BH, Wu W, Cayting P, Boyle AP, Sundaram V, Xing X, Dogan N, Li J, et al.
17 2014. Principles of regulatory information conservation between mouse and human. *Nature* **515**:
18 371–375.

19 Cope L, Zhong X, Garrett E, Parmigiani G. 2004. MergeMaid: R tools for merging and cross-
20 study validation of gene expression data. *Stat Appl Genet Mol Biol* **3**: Article29.

21 De Rie D, Abugessaisa I, Alam T, Arner E, Arner P, Ashoor H, Åström G, Babina M, Bertin N,
22 Burroughs AM, et al. 2017. An integrated expression atlas of miRNAs and their promoters in
23 human and mouse. *Nat Biotechnol* **35**: 872–878.

1 The FANTOM Consortium and the RIKEN PMI and CLST (DGT).2014. A promoter-level
2 mammalian expression atlas. *Nature* **507**: 462–470.

3 Hinrichs AS, Karolchik D, Baertsch R, Barber GP, Bejerano G, Clawson H, Diekhans M, Furey
4 TS, Harte RA, Hsu F, et al. 2006. The UCSC Genome Browser Database: update 2006. *Nucleic
5 Acids Res* **34**(Database issue): D590–D598.

6 Huntley RP, Sawford T, Mutowo-Meullenet P, Shypitsyna A, Bonilla C, Martin MJ, O'Donovan C.
7 2015. The GOA database: Gene Ontology annotation updates for 2015. *Nucleic Acids Res*
8 **43**(Database issue): D1057-D1063.

9 Jubb AW, Young RS, Hume DA, Bickmore WA. 2016. Enhancer turnover is associated with a
10 divergent transcriptional response to glucocorticoid in mouse and human macrophages. *J
11 Immunol* **196**: 813–822.

12 Kin K, Maziarz J, Chavan AR, Kamat M, Vasudevan S, Birt A, Emera D, Lynch VJ, Ott TL,
13 Pavlicev M, Wagner GP. 2016. The transcriptomic evolution of mammalian pregnancy: Gene
14 expression innovations in endometrial stromal fibroblasts. *Genome Biol Evol* **8**: 2459–2473.

15 King MC, Wilson AC. 1975. Evolution at two levels in humans and chimpanzees. *Science* **188**:
16 107–116.

17 Kortenhorst MS, Wissing MD, Rodríguez R, Kachhap SK, Jans JJ, Van der Groep P, Verheul
18 HM, Gupta A, Aiyetan PO, van der Wall E, et al. 2013. Analysis of the genomic response of
19 human prostate cancer cells to histone deacetylase inhibitors. *Epigenetics* **8**: 907–920.

20 Kozomara A, Griffiths-Jones S. 2014. miRBase: annotating high confidence microRNAs using
21 deep sequencing data. *Nucleic Acids Res* **42**(Database issue): D68-D73.

22 Kuhn RM, Haussler D, Kent WJ. 2013. The UCSC genome browser and associated tools. *Brief
23 Bioinform* **14**: 144–161.

24 Lin S, Lin Y, Nery JR, Urich MA, Breschi A, Davis CA, Dobin A, Zaleski C, Beer MA, Chapman

1 WC, et al. 2014. Comparison of the transcriptional landscapes between human and mouse
2 tissues. *Proc Natl Acad Sci U S A* **111**: 17224–17229.

3 Lizio M, Deviatiiarov R, Nagai H, Galan L, Arner E, Itoh M, Lassmann T, Kasukawa T,
4 Hasegawa A, Ros MA, et al. 2017. Systematic analysis of transcription start sites in avian
5 development. *PLoS Biol* **15**: e2002887.

6 Lizio M, Mukarram AK, Ohno M, Watanabe S, Itoh M, Hasegawa A, Lassmann T, Severin J,
7 Harshbarger J, Abugessaisa I, et al. 2017. Monitoring transcription initiation activities in rat and
8 dog. *Sci Data* **4**: 170173.

9 Love MI, Huber W, Anders S. 2014. Moderated estimation of fold change and dispersion for
10 RNA-seq data with DESeq2. *Genome Biol* **15**: 550.

11 Marchionni L, Hayashi M, Guida E, Ooki A, Munari E, Jabboure FJ, Dinalankara W, Raza A,
12 Netto GJ, Hoque MO, Argani P. 2017. MicroRNA expression profiling of Xp11 renal cell
13 carcinoma. *Hum Pathol* **67**: 18–29.

14 Merkin J, Russell C, Chen P, Burge CB. 2012. Evolutionary dynamics of gene and isoform
15 regulation in mammalian tissues. *Science* **338**: 1593–1599.

16 Musser JM, Wagner GP. 2015. Character trees from transcriptome data: Origin and
17 individuation of morphological characters and the so-called "species signal". *J Exp Zool B Mol*
18 *Dev Evol* **324**: 588–604.

19 NCBI Resource Coordinators 2018. Database resources of the National Center for
20 Biotechnology Information. *Nucleic Acids Res* **46**(D1): D8–D13.

21 Odom DT, Dowell RD, Jacobsen ES, Gordon W, Danford TW, MacIsaac KD, Rolfe PA, Conboy
22 CM, Gifford DK, Fraenkel E. 2007. Tissue-specific transcriptional regulation has diverged
23 significantly between human and mouse. *Nat Genet* **39**: 730–732.

24 Pachkov M, Balwierz PJ, Arnold P, Ozonov E, Van Nimwegen E. 2013. SwissRegulon, a

1 database of genome-wide annotations of regulatory sites: recent updates. *Nucleic Acids Res*
 2 **41**(Database issue): D214–D220.

3 Parmigiani G, Garrett-Mayer ES, Anbazhagan R, Gabrielson E. 2004. A cross-study comparison
 4 of gene expression studies for the molecular classification of lung cancer. *Clin Cancer Res* **10**:
 5 2922–2927.

6 Pishesha N, Thiru P, Shi J, Eng JC, Sankaran VG, Lodish HF. 2014. Transcriptional divergence
 7 and conservation of human and mouse erythropoiesis. *Proc Natl Acad Sci U S A* **111**: 4103–
 8 4108.

9 Ramilowski JA, Goldberg T, Harshbarger J, Kloppmann E, Lizio M, Satagopam VP, Itoh M,
 10 Kawaji H, Carninci P, Rost B, Forrest AR. 2015. A draft network of ligand-receptor-mediated
 11 multicellular signalling in human. *Nat Commun* **6**: 7866.

12 Regev A, Teichmann SA, Lander ES, Amit I, Benoist C, Birney E, Bodenmiller B, Campbell P,
 13 Carninci P, Clatworthy M, et al. 2017. The human cell atlas. *eLife* **6**: e27041.

14 Ross AE, Marchionni L, Phillips TM, Miller RM, Hurley PJ, Simons BW, Salmasi AH, Schaeffer
 15 AJ, Gearhart JP, Schaeffer EM. 2011. Molecular effects of genistein on male urethral
 16 development. *J Urol* **185**: 1894–1898.

17 Schroder K, Irvine KM, Taylor MS, Bokil NJ, Le Cao KA, Masterman KA, Labzin LI, Semple CA,
 18 Kapetanovic R, Fairbairn L, et al. 2012. Conservation and divergence in Toll-like receptor 4-
 19 regulated gene expression in primary human versus mouse macrophages. *Proc Natl Acad Sci U*
 20 *S A* **109**: E944-E953.

21 Su AI, Cooke MP, Ching KA, Hakak Y, Walker JR, Wiltshire T, Orth AP, Vega RG, Sapinoso
 22 LM, Moqrich A, et al. 2002. Large-scale analysis of the human and mouse transcriptomes. *Proc*
 23 *Natl Acad Sci U S A* **99**: 4465–4470.

24 Suzuki H, Forrest ARR, Van Nimwegen E, Daub CO, Balwierz PJ, Irvine KM, Lassmann T,

1 Ravasi T, Hasegawa Y, De Hoon MJL, et al. 2009. The transcriptional network that controls
2 growth arrest and differentiation in a human myeloid leukemia cell line. *Nat Genet* **41**: 553–562.

3 The ENCODE Project Consortium. 2012. An integrated encyclopedia of DNA elements in the
4 human genome. *Nature* **489**: 57–74.

5 Tyekucheva S, Marchionni L, Karchin R, Parmigiani G. 2011. Integrating diverse genomic data
6 using gene sets. *Genome Biol* **12**: R105.

7 Vickaryous MK, Hall BK. 2006. Human cell type diversity, evolution, development, and
8 classification with special reference to cells derived from the neural crest. *Biol Rev Camb Philos*
9 *Soc* **81**: 425–455.

10 Vilella AJ, Severin J, Ureta-Vidal A, Heng L, Durbin R, Birney E. 2009. EnsemblCompara
11 GeneTrees: Complete, duplication-aware phylogenetic trees in vertebrates. *Genome Res* **19**:
12 327–335.

13 Yang F, Chen Q, He S, Yang M, Maguire EM, An W, Afzal TA, Luong LA, Zhang L, Xiao Q.
14 2018. miR-22 is a novel mediator of vascular smooth muscle cell phenotypic modulation and
15 neointima formation. *Circulation* **137**: 1824–1841.

16 Young RS, Hayashizaki Y, Andersson R, Sandelin A, Kawaji H, Itoh M, Lassmann T, Carninci P,
17 FANTOM Consortium, Bickmore WA, et al. 2015. The frequent evolutionary birth and death of
18 functional promoters in mouse and human. *Genome Res* **25**: 1546–1557.

19 Yue F, Cheng Y, Breschi A, Vierstra J, Wu W, Ryba T, Sandstrom R, Ma Z, Davis C, Pope BD,
20 et al. 2014. A comparative encyclopedia of DNA elements in the mouse genome. *Nature* **515**:
21 355–364.

22 Zerbino DR, Achuthan P, Akanni W, Amode MR, Barrell D, Bhai J, Billis K, Cummins C, Gall A,
23 Girón CG, et al. 2018. Ensembl 2018. *Nucleic Acids Res* **46**(D1): D754–D761.

1 Figure captions

2 **Figure 1.** Gene expression PCA. (A) PCA for all samples of cell types in common between
3 human and mouse. (B) PCA for all samples of cell types in common between all five species.

4

5 **Figure 2.** Differential gene expression analysis. (A) Expression profile of *HNF4A*, *ELF2*, and
6 *FOXO1* as examples of genes with an expression profile highly correlated across cell types
7 between species, but with significant expression level differences between species in specific
8 cell types. (B) Cumulative distribution of Pearson's correlation r across cell types in gene
9 expression between human and mouse, rat, dog, or chicken. The number N of expressed
10 orthologous genes included in the distribution is shown in the vertical axis label, and the
11 estimated median value of r is indicated on the horizontal axis of each graph. The background
12 distribution of r obtained by randomizing genes (solid) or randomizing samples (dashed) as well
13 as the theoretical background distribution of r for an uncorrelated bivariate normal distribution
14 (dotted) are shown in grey. The statistical significance was calculated using the Mann-Whitney
15 U test comparing Pearson's correlation values for orthologs to the background distribution of r
16 for randomly paired genes between human and mouse, rat, dog, or chicken. Note that the
17 median correlation values are not directly comparable between species, as the sets of
18 orthologous genes are different. (C) Differential gene expression analysis of orthologous genes
19 in human compared to mouse, rat, dog, and chicken. The red and blue bars correspond to the
20 percentage of expressed orthologous genes with significantly (Benjamini-Hochberg corrected p
21 < 0.1) higher and lower expression, respectively, in human compared to mouse, rat, dog, or
22 chicken. The number N of orthologous genes expressed in each cell type is shown on the right.

23

Figure 3. Promoter analysis of differentially expressed genes. (A) Percentage of genes in mouse, rat, dog, and chicken for which the dominant promoter was located in a genome region that had an orthologous genome region in human, and the percentage that the orthologous region contained the dominant promoter for the orthologous gene in human. (B) Percentage of differentially expressed genes in each cell type depending on whether the genomic region of the dominant promoter in each species had an orthologous genomic region in the human genome. The one-sided p -value calculated using Fisher's exact test is shown on the right, together with the number N of expressed genes in each cell type.

9

Figure 4. Conservation and Gene Ontology analysis of differentially expressed genes. (A) Percentage of differentially expressed genes in each cell type as a function of age of the most recent common ancestor. The one-sided p -value of a Poisson regression model against the evolutionary age category is shown on the right, together with the number N of expressed genes in each cell type with an annotation in the NCBI HomoloGene database. (B) Gene Ontology analysis of differentially expressed genes. The p -value, calculated using Fisher's exact test, of overrepresentation or underrepresentation of differentially expressed genes in each Gene Ontology category compared to an expression-matched set of background genes is shown in red and blue, respectively.

19

Figure 5. Motif Activity analysis. (A-D) Examples of calculated motif activities in human and mouse for motifs associated with the broadly expressed transcription factor TP53 (A), the hematopoietic lineage-specific RUNX transcription factors (B), the hepatocyte-specific HNF4A transcription factor (C), and the testis-specific transcription factor SPZ1 (D). Each of the 15 matching cell types between human and mouse is shown as a dot. The blood cell types CD19⁺ B cells, CD4⁺ T cells, CD8⁺ T cells, common myeloid progenitors, and granulocyte macrophage progenitors are shown in red for the RUNX motif, and the liver cell types hepatic sinusoidal

1 endothelial cells, hepatic stellate cells (lipocytes), and hepatocytes are shown in green for the
 2 motif associated with HNF4A. (E-F) Cumulative distribution of Pearson's correlation r across cell
 3 types in motif activity for promoters (E) and enhancers (F) between human and mouse, rat, dog,
 4 and chicken. The estimated median value of r is indicated on the horizontal axis of each graph.
 5 As a background distribution, we calculated the same correlation between pairs of different
 6 motifs in human and mouse, rat, dog, and chicken. The Mann-Whitney U test p -value comparing
 7 the actual correlation values to the correlation values of the background distribution is shown for
 8 each comparison.

9

10 **Figure 6.** Differential miRNA expression analysis. (A) Differential expression analysis of
 11 miRNAs using FANTOM5 sRNA sequencing data in aortic smooth muscle cells in human
 12 compared to mouse, rat, dog or chicken. The red and blue bars correspond to the percentage of
 13 expressed orthologous miRNAs with significantly (Benjamini-Hochberg corrected $p < 0.1$) higher
 14 and lower expression, respectively, in human compared to mouse, rat, dog, or chicken. The
 15 number N of expressed orthologous miRNAs in each comparison is shown on the right. (B)
 16 Differential expression analysis of miRNAs in human compared to mouse, rat, dog, and chicken;
 17 using CAGE expression of the pri-miRNA as a proxy for the expression level of the mature
 18 miRNA. The red and blue bars correspond to the percentage of expressed orthologous miRNAs
 19 with significantly (Benjamini-Hochberg corrected $p < 0.1$) higher and lower expression,
 20 respectively, in human compared to mouse, rat, dog, or chicken. The number N of expressed
 21 orthologous miRNAs in each comparison is shown on the right. (C) Percentage of miRNAs
 22 differentially expressed in each comparison, separately based on the evolutionary age of each
 23 miRNA. The one-sided p -value of a Poisson regression model against the evolutionary age
 24 category is shown on the right, together with the number N of expressed orthologous miRNAs in
 25 each comparison. (D) Percentage of miRNAs differentially expressed in each comparison,
 26 separately based on the evolutionary age of each miRNA; using CAGE expression of the pri-

1

- 1 miRNA as a proxy for the expression level of the mature miRNA. The one-sided p -value of a
- 2 Poisson regression model against the evolutionary age category is shown on the right, together
- 3 with the number N of expressed orthologous miRNAs in each comparison.

1 Supplemental Methods

2 Gene expression quantitation

3 For each promoter region for each gene, we found the total CAGE tag count, summing over all
 4 libraries included in this study. For each gene, the region with the highest total CAGE tag count
 5 was identified as the dominant promoter, and the position with the highest CAGE tag count
 6 within this region as the representative transcription start site of the gene. For each gene, the
 7 CAGE expression level in a CAGE library was defined as the sum of CAGE tags across all
 8 regions associated with the gene. CAGE expression data were normalized to tags-per-million
 9 (t.p.m.) by dividing by the total sum of CAGE tags associated with genes, and multiplying by
 10 1,000,000. Genes with an expression level of at least 10 t.p.m. were included in the analysis.

11 RNA-seq expression data

12 ENCODE (The ENCODE Project Consortium 2012) RNA-seq gene expression data were
 13 downloaded from the ENCODE portal (Davis et al. 2018) (<https://www.encodeproject.org/>) as
 14 gene quantifications in tsv format with sample annotations in json format for the following
 15 accession numbers: ENCFF005MLW, ENCFF051UVH, ENCFF084WWG, ENCFF115BZB,
 16 ENCFF122LQH, ENCFF134RPA, ENCFF143IYG, ENCFF152DMV, ENCFF162KBI,
 17 ENCFF201BIE, ENCFF203RCU, ENCFF227RBV, ENCFF280HQA, ENCFF362GHJ,
 18 ENCFF375LAX, ENCFF421SLK, ENCFF428HNN, ENCFF442XPO, ENCFF456GZP,
 19 ENCFF485EYP, ENCFF507RNA, ENCFF508QDF, ENCFF509LQW, ENCFF517ZTC,
 20 ENCFF547TAC, ENCFF549OGG, ENCFF603OKL, ENCFF625HJC, ENCFF637YVK,
 21 ENCFF664JKG, ENCFF677ZIY, ENCFF686JQP, ENCFF750HMK, ENCFF780DAF,
 22 ENCFF784ZTQ, ENCFF798PSE, ENCFF804BIT, ENCFF804WTK, ENCFF809RAX,

1

1 ENCFF850LMK, ENCFF866LBS, ENCFF878LJT, ENCFF878UHQ, ENCFF908GIP,
2 ENCFF911YUO, ENCFF916CFV, ENCFF926FDN, ENCFF937HOV. RNA-seq gene expression
3 count tables for endometrial stromal fibroblast (Kin et al. 2016), their normalized expression
4 levels in tags per million, and gene ortholog associations were downloaded from the NCBI Gene
5 Expression Omnibus under accession number GSE67659.

6 Pairwise genome alignments

7 Preprocessing of the genome sequences and postprocessing of the alignments was performed
8 using partitionSequence.pl, blastz-normalizeLav, lavToPsl, axtChain, chainAntiRepeat,
9 chainMergeSort, chainPreNet, chainNet, netSyntenic, netFilter, netToAxt, axtSort, and
10 axtToMaf, which are part of the UCSC Genome Browser bioinformatics utilities (Kuhn et al.
11 2013) release 366 (June 5, 2018). First, for each pairwise alignment we used
12 partitionSequence.pl to split the target and query genome sequences into segments of 20 Mb,
13 with the target sequence segments overlapping by 10 kb. Sequences were aligned using lastz
14 (Harris 2007) version 1.03.66 with parameters shown in Supplemental Table S15. Alignment
15 coordinates were corrected using blastz-normalizeLav and converted to .psl format using
16 lavToPsl. Alignments were chained using axtChain with parameters shown in Supplemental
17 Table S15 and further processed using chainAntiRepeat, chainMergeSort, chainPreNet,
18 chainNet, netSyntenic, and netFilter. The best alignments were extracted using netToAxt,
19 followed by axtSort and axtToMaf to generate a .maf (multiple alignment format) file.

20 Genome-wide TFBS prediction

21 For human and mouse, alignments for human, macaque, mouse, rat, dog, horse, cow,
22 opossum, and chicken were extracted from the 100-way (human) and 30-way (mouse, rat, dog,
23 and chicken) multiple genome alignments; for rat, dog, and chicken, we used the 30-way

1

1 multiple genome alignments described in the main Methods section of this manuscript.
2 Extracted alignments were improved as described previously (Arner et al. 2015) using the T-
3 Coffee (Notredame et al. 2000) multiple sequence aligner version 9.01. We ran MotEvo (Arnold
4 et al. 2012) version 1.01 on these alignments as described previously (Arner et al. 2015), but
5 using the 190 motifs in the SwissRegulon (Pachkov et al. 2013) release of July 13, 2015, and
6 with a background probability for the four nucleotides that was calculated by counting their
7 frequency across the genome.

8 Motif activity analysis

9 The scripts `make_profile.py`, `associate_tfbs.py`, and `calculate_motif_activity.py`, available at
10 http://fantom.gsc.riken.jp/5/suppl/Alam_et_al_2020/ and in the Supplemental Code, provide an
11 implementation of the motif activity analysis methodology described previously (Suzuki et al.
12 2009). The density of predicted transcription factor binding sites with respect to the
13 representative transcription start site was calculated for each motif was calculated using the
14 script `make_profile.py` with default options. Using the calculated density, predicted binding sites
15 were associated with the dominant promoter of each gene by the script `associate_tfbs.py` with
16 default options. The density of predicted transcription factor binding sites in a ± 500 base pair
17 window around enhancers was calculated using the script `make_profile.py` with options `--`
18 `upstream=500 --downstream=500 --symmetric`. Using the calculated density, predicted binding
19 sites were associated with each enhancer by the script `associate_tfbs.py` with default options.
20 We then selected motifs with at least 50 predicted binding sites both for promoters and in
21 enhancers in each pair of species, resulting in 154 (human-mouse), 151 (human-rat), 149
22 (human-dog), and 94 (human-chicken) motifs. Motif activities were calculated for promoters and
23 enhancers separately for the selected motifs using the script `calculate_motif_activity.py` with
24 option `'-n 0'`.

1 MicroRNA analysis

2 We used release 21 of the miRBase database (Kozomara and Griffiths-Jones 2014), lifted over
3 for rat and chicken to genome assembly rn6 and galGal5 respectively, as our reference set of
4 known miRNAs; one pre-miRNA in rat and 26 pre-miRNAs in chicken could not be lifted over
5 and were dropped. Pre-miRNAs were classified as robust or permissive (Supplemental Table
6 S9) as described previously (De Rie et al. 2017).

7 Candidate novel miRNAs were identified using miRDeep2 (Friedländer et al. 2012), resulting in
8 229 (rat), 249 (dog), and 180 (chicken) predicted pre-miRNAs, including 169 (rat), 179 (dog),
9 and 128 (chicken) known pre-miRNAs and 59 (rat), 74 (dog), and 55 (chicken) candidate novel
10 pre-miRNAs (Supplemental Table S10).

11 Promoters for miRNAs in rat, dog, and chicken were identified using the same approach as
12 applied previously for human and mouse (De Rie et al. 2017), using transcripts annotated for
13 rat, dog, or chicken in the NCBI Entrez Gene database (Brown et al. 2015), downloaded on 13
14 November 2017, as candidate primary miRNAs. Identified promoters for all robust miRNAs in
15 mouse, rat, dog, and chicken were curated manually by two annotators (Supplemental Table
16 S12).

17 Supplemental References

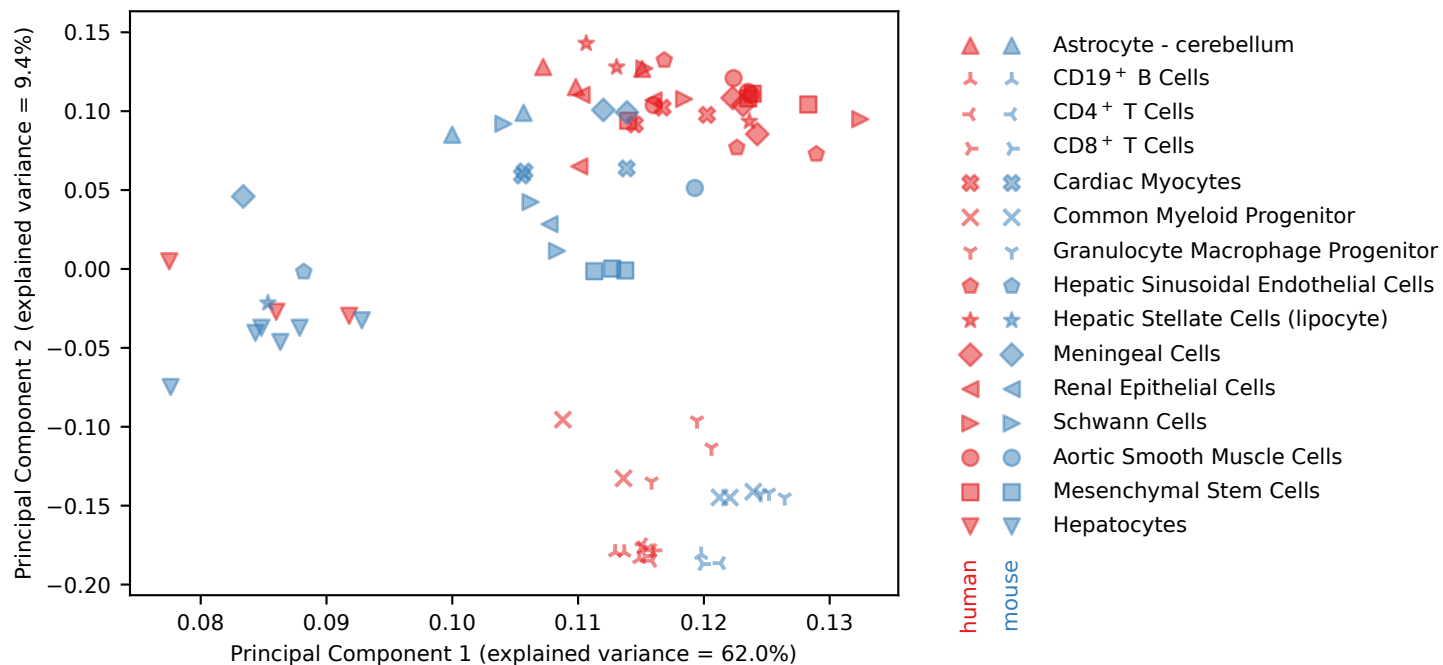
18 Davis CA, Hitz BC, Sloan CA, Chan ET, Davidson JM, Gabdank I, Hilton JA, Jain K,
19 Baymuradov UK, Narayanan AK, et al. 2018. The Encyclopedia of DNA elements (ENCODE):
20 data portal update. *Nucleic Acids Res* **46**(D1): D794–D801.

21 Friedländer MR, Mackowiak SD, Li N, Chen W, Rajewsky N. 2012. miRDeep2 accurately
22 identifies known and hundreds of novel microRNA genes in seven animal clades. *Nucleic Acids*
23 *Res* **40**: 37–52.

- 1 Harris RS. 2007. Improved pairwise alignment of genomic DNA. Ph.D. thesis. Pennsylvania
- 2 State University, State College, PA.
- 3 Notredame C, Higgins DG, Heringa J. 2000. T-Coffee: A novel method for fast and accurate
- 4 multiple sequence alignment. *J Mol Biol* **302**: 205–217.

Figure 1

A



B

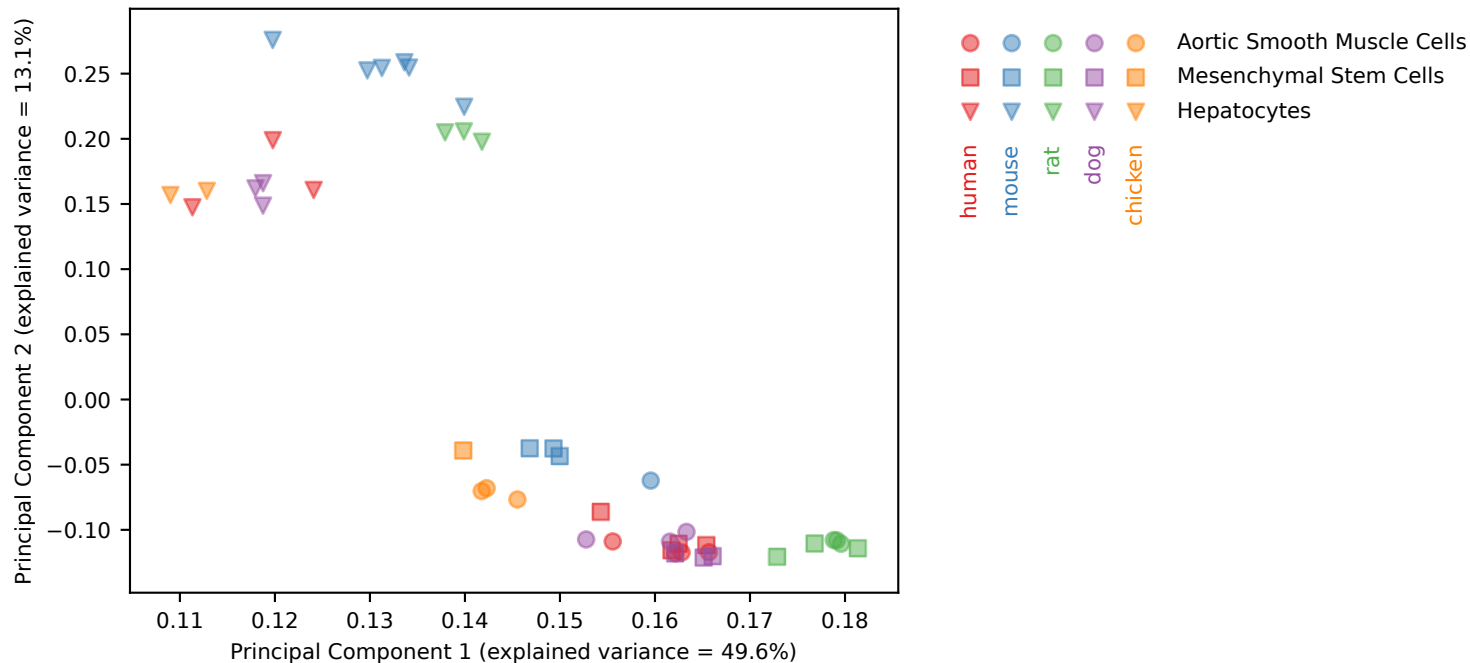


Figure 2

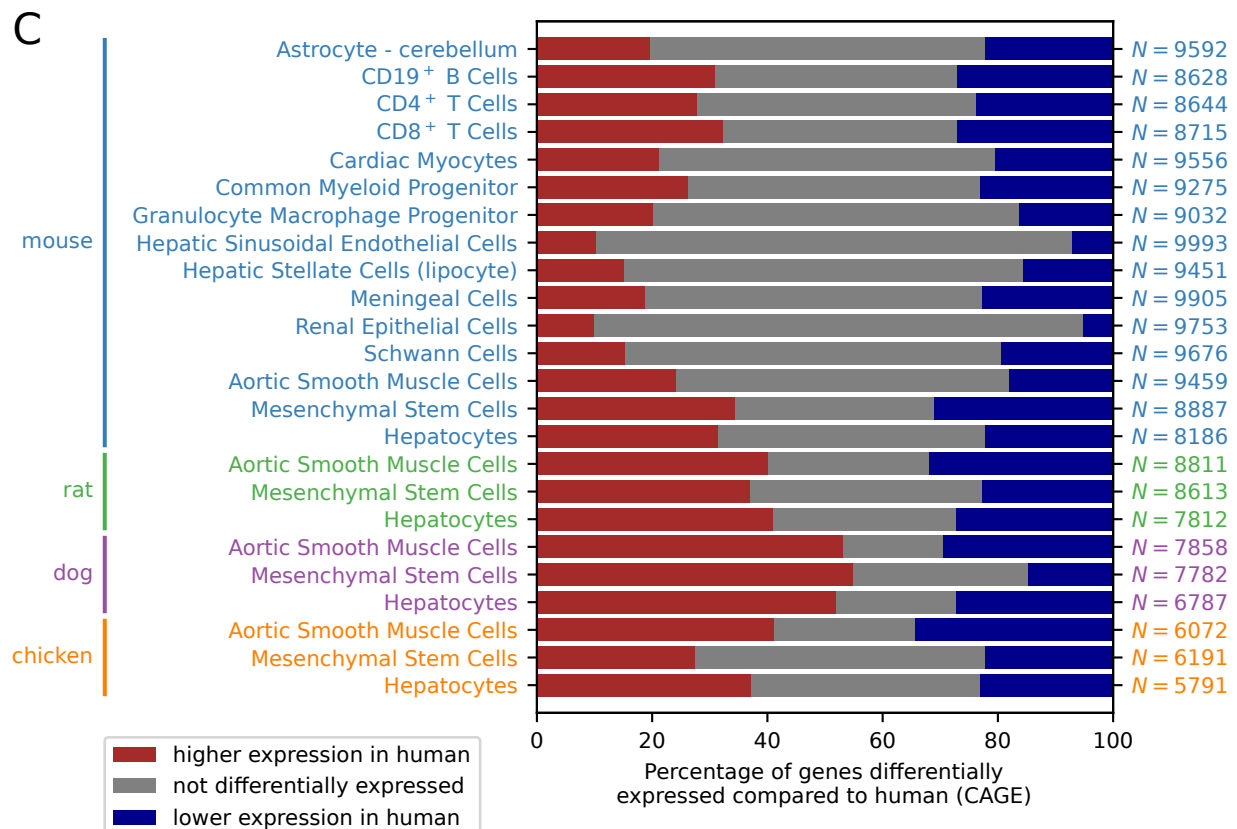
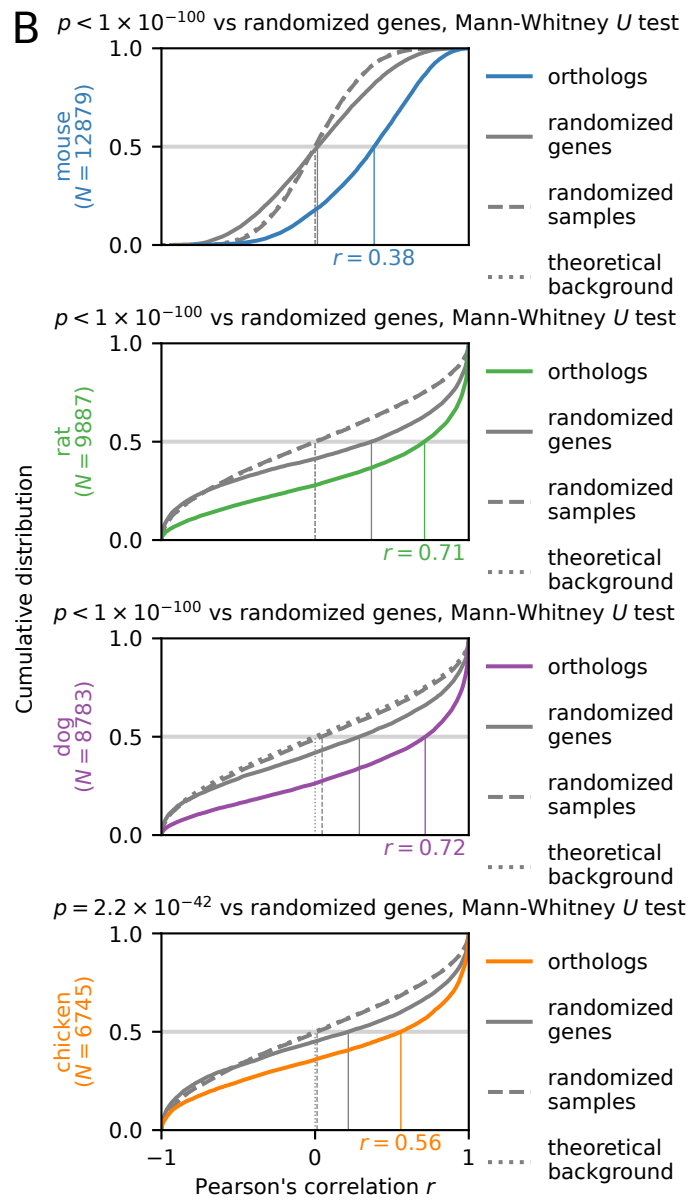
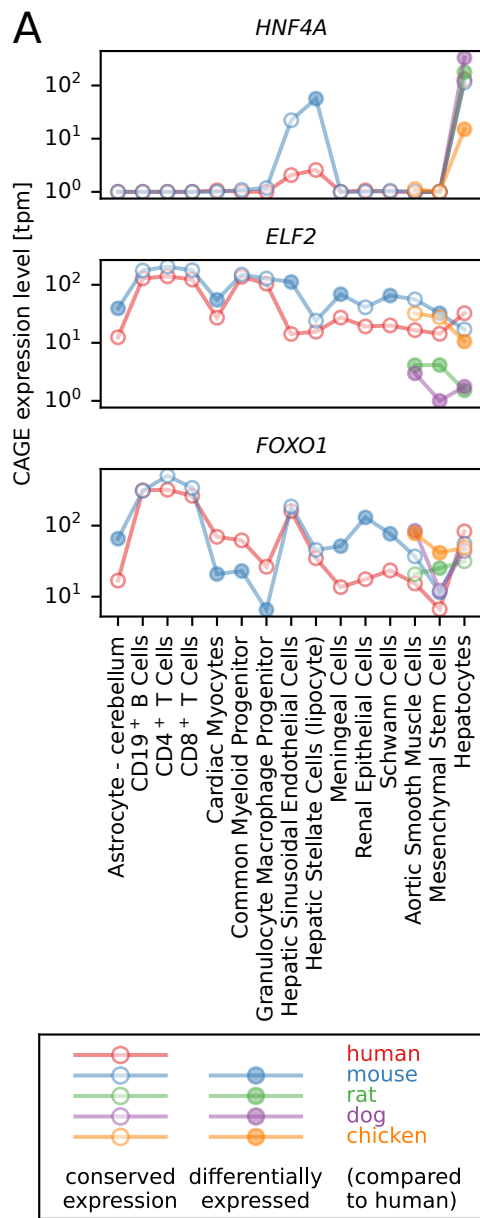


Figure 3

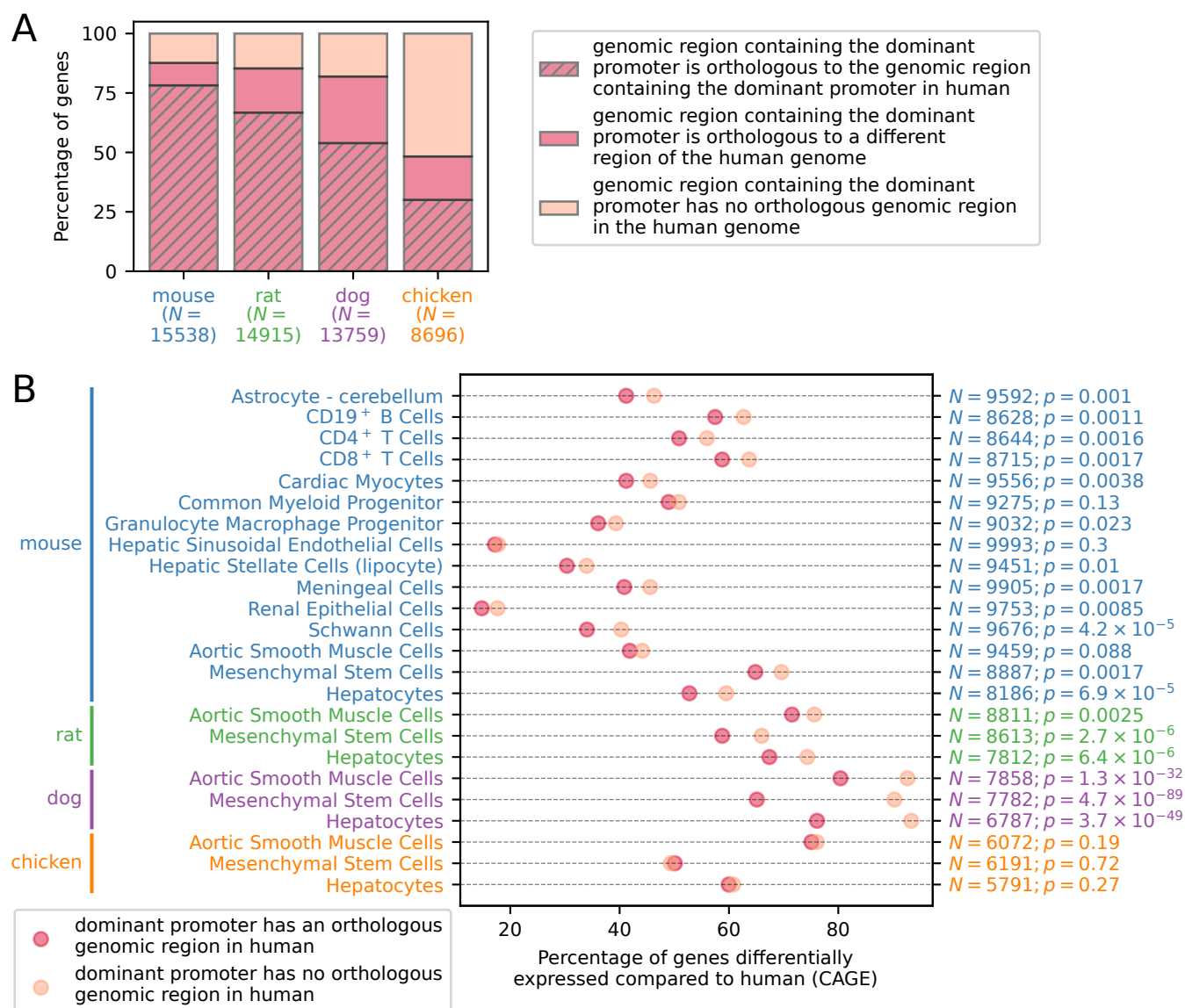


Figure 4

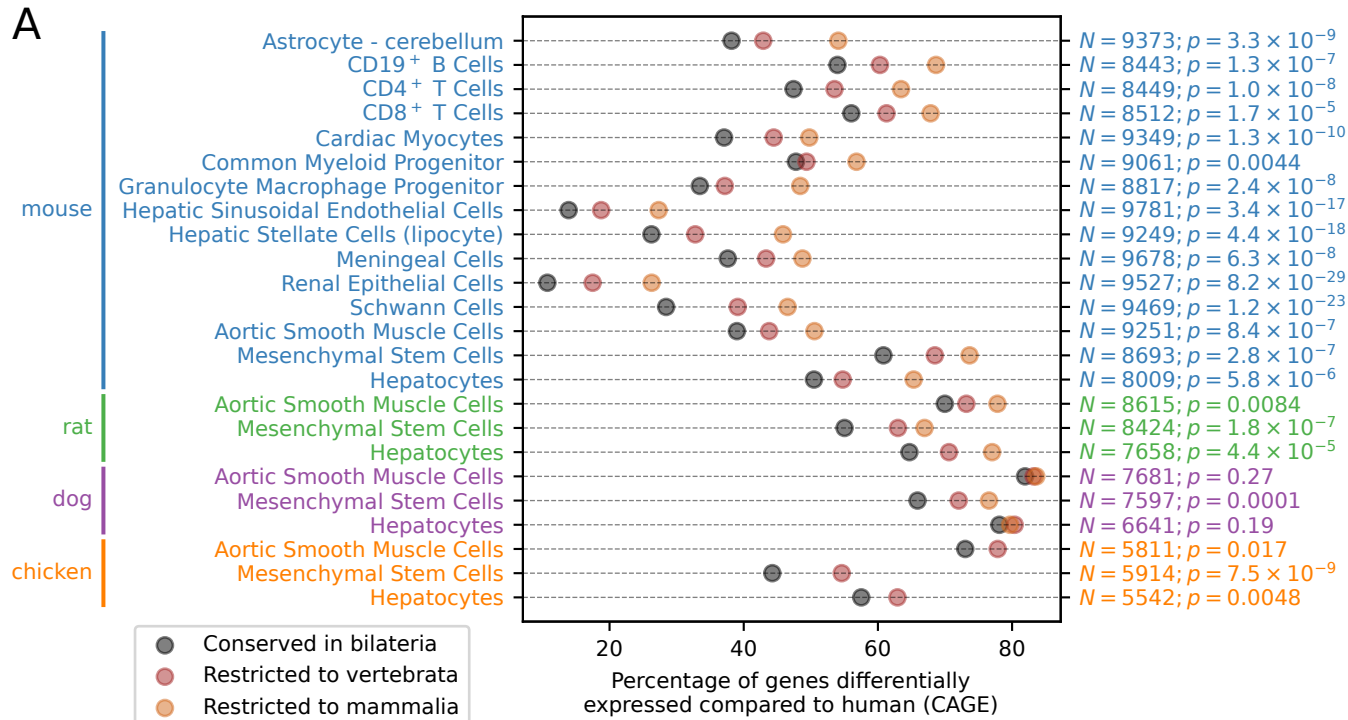
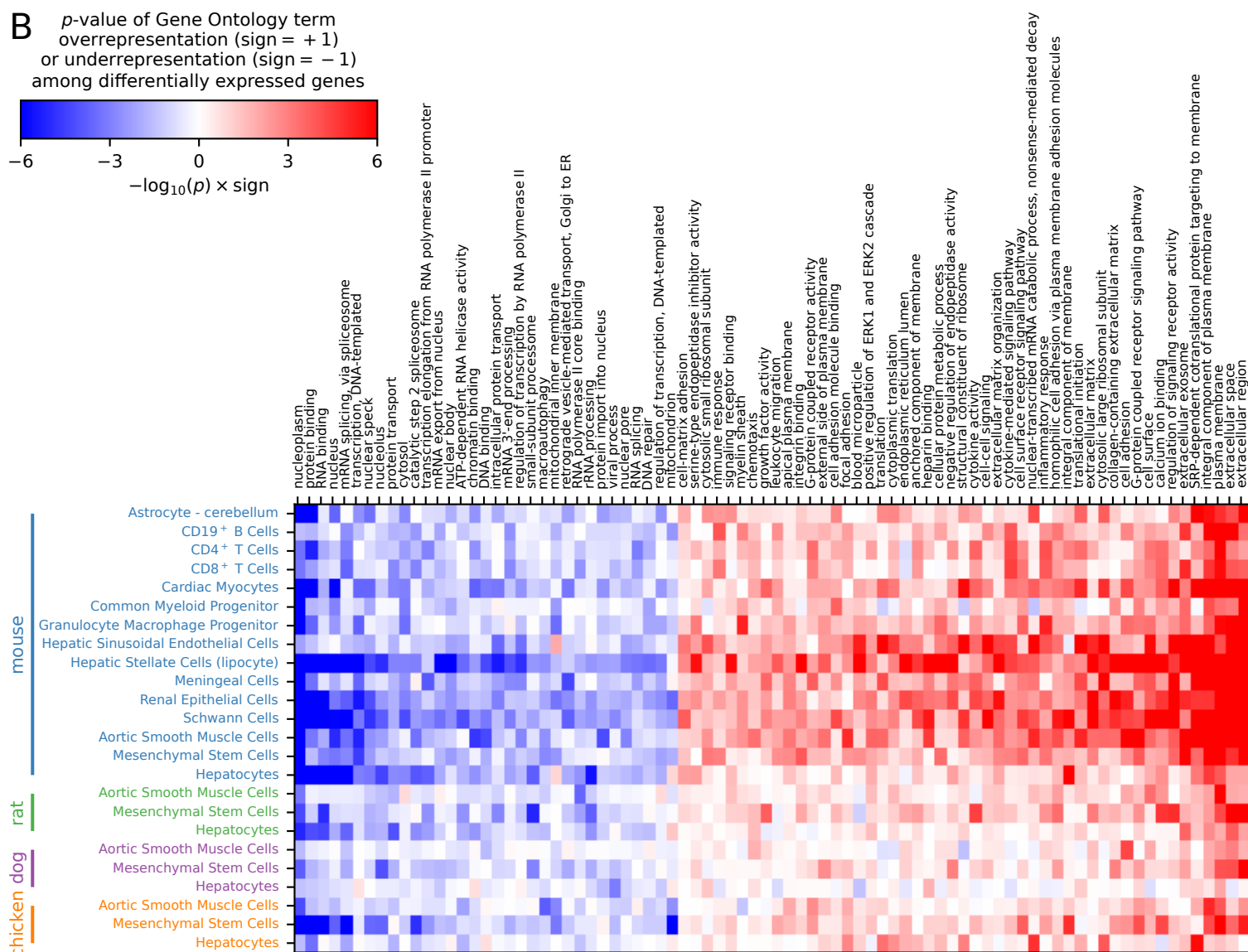
A**B**

Figure 5

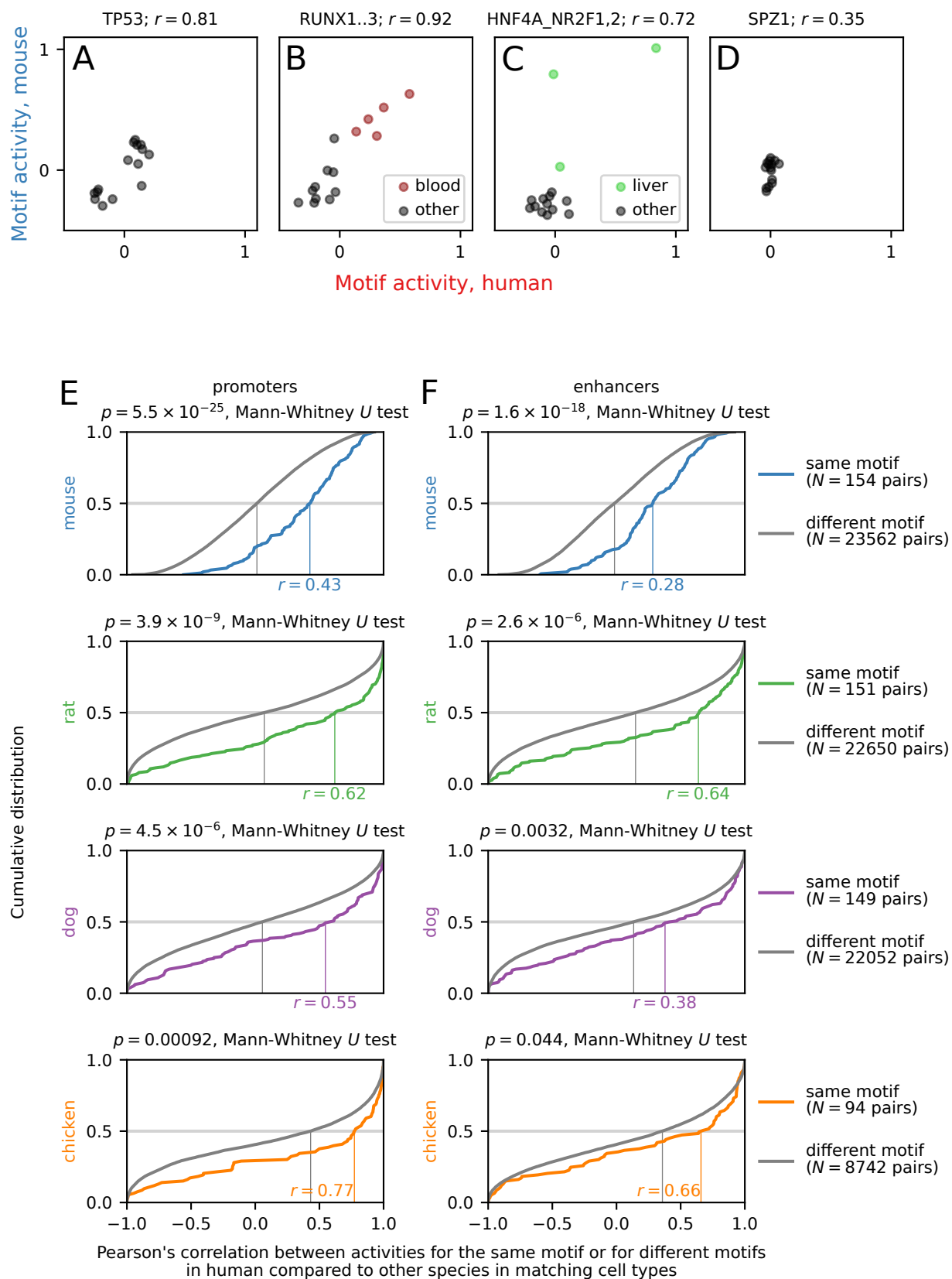


Figure 6

